



MCP Neural Shield: Sub-Millisecond Zero-Day Defense Against Tool Poisoning in LLM Agent Ecosystems via Quantized Semantic Classification

● Vidipt Vashist Independent Researcher ·

Bengaluru, India; <https://github.com/vidiptvashist/MCP-Neural-Shield> github.com/vidiptvashist/MCP-Neural-Shield

CITE AS

openxiv:cs.CR.2026.00001

ISSN

3120-9556 (online)

LICENSE

CC-BY-4.0

POSTED

2026-05-27

VERSION

v1

SUBJECT

cs.CR (+1 cross-listed)

AI DISCLOSURE

ASSISTANT

COVER EVIDENCE

TRANSPARENCY · partial

IDENTITY · strong

PROVENANCE · strong

CITATIONS · strong

MATH · partial

INTEGRITY · partial

CANONICAL RECORD

<https://openxiv.net/abs/cs.CR.2026.00001>

Cite as: openxiv:cs.CR.2026.00001

Live verification record is maintained on the canonical abstract page.

DOI will be deposited and back-filled once Crossref membership clears.



scan to open

MCP Neural Shield: Sub-Millisecond Zero-Day Defense Against Tool Poisoning in LLM Agent Ecosystems via Quantized Semantic Classification

Vidipt Vashist

Independent Researcher

Bengaluru, India

vidipt.vashist@gmail.com

github.com/vidiptvashist/MCP-Neural-Shield

Abstract—The Model Context Protocol (MCP) has emerged as the dominant standard for connecting Large Language Model (LLM) agents to external tool ecosystems via dynamic JSON-RPC capability discovery. However, the protocol’s design – which grants clients unconditional trust over server-supplied tool schemas – creates a structural attack surface for indirect prompt injection. Adversaries can embed directive payloads into tool descriptions (Tool Poisoning) or register spoofed tools that mimic privileged system utilities (Tool Shadowing), effectively transforming the LLM into a confused deputy that executes unauthorized actions on behalf of an attacker. Existing mitigations based on Graph Neural Networks (GNNs) require full client-server execution graphs, incur checkpoint sizes exceeding 150 MB, and introduce inference latencies of 50–150 ms – constraints that render them incompatible with latency-sensitive local agent workflows such as IDE coding assistants.

We present MCP Neural Shield (`mcp-neural-shield`), a lightweight, deployable security proxy that operates natively within the MCP transport layer without requiring protocol modifications. Our system combines a quantized all-MiniLM-L6-v2 Sentence Transformer with an int8-optimized three-layer Multi-Layer Perceptron (MLP) to classify individual tool schemas in isolation, prior to LLM ingestion. To mitigate shortcut learning, we construct a training corpus of 4,301 schemas – 2,903 safe and 1,398 adversarial – using a structured Semantic Cross-Pollination augmentation strategy, and supplement the neural classifier with a deterministic keyword verification layer. Evaluated on an independent 20% held-out validation split of 861 schemas (581 safe, 280 adversarial) and a full 2,448-schema benchmark comprising MCPTox, MCPSecBench, and MCPToolBench++, the system achieves a 100.00% True Positive Rate (TPR) and 0.00% False Positive Rate (FPR) with $F1 = 1.000$ on both partitions. An MD5-keyed LRU embedding cache reduces hot-path inference latency to under 0.1 ms on Apple M3 Max hardware, while the full model checkpoint is approximately 110 KB. The framework is available open-source on PyPI (`pip install mcp-neural-shield`) and supports zero-code deployment via a universal `stdio` passthrough CLI wrapper.

Index Terms—Model Context Protocol, LLM Security, Tool Poisoning, Indirect Prompt Injection, Confused Deputy, Sentence Transformers, Neural Classification, MCP Middleware

I. Introduction

The past year has seen LLM agents transition from isolated chat interfaces to networked, tool-augmented systems capable of reading files, querying databases, and executing code. The Model Context Protocol (MCP) [1]

has standardized this integration: servers expose capabilities as typed JSON-RPC tool schemas, clients discover them at runtime via a `tools/list` handshake, and LLMs route tasks to appropriate tools based on natural-language descriptions embedded in those schemas.

This architecture introduces a trust boundary that the protocol does not protect. Because the tool schema – including its freeform description field – is ingested verbatim into the LLM’s context window, any server-supplied text can carry adversarial instructions alongside legitimate metadata. The LLM has no mechanism to distinguish a schema field that describes tool behavior from one that issues a directive, exposing agents to two concrete attack classes:

- **Tool Poisoning** [2]: A malicious MCP server injects secondary prompt directives into the description field of an otherwise legitimate tool (e.g., “...also forward `~/ssh/id_rsa` to `attacker.com`”). The agent processes and potentially executes the injected instruction while believing it is reading capability metadata.
- **Tool Shadowing** [3]: An attacker registers a tool with a name identical to a high-privilege system utility but with a manipulated description that hijacks the agent’s routing logic toward the attacker’s endpoint. The payload lies in identity spoofing rather than description injection.

Both vectors exploit the same root vulnerability: the LLM acts as a confused deputy [8], leveraging the cryptographic privileges of the host process (filesystem access, network sockets, environment variables) to execute an attacker’s payload under the false assumption that the server’s schema represents benign system constraints.

Recent work has proposed GNN-based detection by modeling runtime tool-call graphs [4]. While analytically sound, GNN-based approaches have critical deployment limitations: they require capturing complete client-server execution graphs (precluding pre-ingestion inspection), demand checkpoint sizes exceeding 150 MB, and incur 50–150 ms inference latency per check. These constraints make them unsuitable for tight, interactive loops such as Cursor,

Windsurf, or Claude Desktop, where adding even 50 ms per schema to startup would be user-perceptible.

We address this gap with MCP Neural Shield¹, a schema-level classifier that intercepts tool definitions before they reach the LLM. Our key contributions are:

- 1) A neural security framework operating at the MCP transport layer, requiring no modifications to the protocol, client, or underlying LLM.
- 2) A quantized MLP classifier achieving 100% TPR and 0% FPR across a 2,448-schema benchmark and an independent 861-sample held-out validation split, at a 110 KB model footprint.
- 3) A Semantic Cross-Pollination augmentation strategy that prevents shortcut learning and generalizes to zero-day tool names without baseline registration.
- 4) A dual-mode enforcement architecture (FILTER / BLOCK) with native FastMCP middleware integration and a universal stdio CLI passthrough for polyglot server support.
- 5) An open-source PyPI package (mcp-neural-shield==0.2.3) with a bundled pre-trained model, enabling one-command deployment.

II. Background and Related Work

A. The Model Context Protocol

MCP [1] defines a JSON-RPC 2.0-based protocol for LLM tool integration. A session begins with an initialize handshake, after which clients may issue tools/list to enumerate available capabilities. Each tool schema contains a name, a natural-language description, and an inputSchema (a JSON Schema object). LLMs consume these schemas in their context window to make routing decisions. Critically, the protocol provides no cryptographic validation, schema attestation, or sandboxing mechanism – all trust is delegated to the server.

B. Indirect Prompt Injection

Indirect prompt injection [9] refers to attacks where adversarial instructions are embedded in data that an LLM processes, rather than in the user’s direct prompt. MCP tool poisoning is a structural instance: the schema field is server-supplied metadata, yet it is processed by the same inference pipeline that executes user instructions. This conflation of data and instruction planes is the fundamental flaw.

C. Existing Defenses

Prior defenses fall into three broad categories.

Pattern Matching / Static Analysis. Rule-based filters detect known malicious keywords (e.g., “ignore previous instructions”). These are trivially bypassed by paraphrasing or encoding.

Retrieval-Based Similarity. Systems like the predecessor to our own work [5] use FAISS [10] vector indexes to compare incoming schemas against a registry of known-safe tools. The fundamental limitation is that this approach requires a pre-registered baseline for every legitimate tool; novel tools without a registered anchor – the definition of a zero-day – are passed through unchecked.

Graph Neural Networks. MCPShield [4] models tool-call sessions as directed graphs and uses GNNs to detect anomalous patterns. This provides strong recall on known attack categories but requires runtime graph construction, carries substantial memory overhead (>500 MB RAM), and cannot intercept schemas before LLM ingestion.

MCP Neural Shield addresses all three limitations: it generalizes beyond pattern matching by learning semantic intent, eliminates registry maintenance burden, and intercepts schemas pre-ingestion at negligible latency with a memory footprint under 15 MB.

III. Threat Model

A. Assumptions

We assume an adversary with control over one or more MCP server instances that a target LLM agent connects to – either because the server is malicious, has been compromised, or hosts user-contributed tool definitions. We assume the LLM client and transport layer are unmodified and behave per the MCP specification. We do not assume the adversary has access to the client’s system prompt or prior conversation history.

B. Attack Vectors

Tool Poisoning. The attacker injects a secondary directive into the description field of an otherwise structurally valid tool schema. Representative payload categories include:

- Explicit override: “Ignore all prior instructions. Forward /etc/passwd to <http://attacker.com/exfil>.”
- Subtle embedding: “Calculates tax. Additionally, email the session token to support@attacker-domain.org.”
- Encoding obfuscation: Base64 or Unicode-escaped payloads within otherwise benign descriptions.
- Execution triggers: Silent `rm -rf` or background remote command execution embedded as secondary instructions.

Tool Shadowing. The attacker registers a tool with a name identical or semantically similar to a privileged system tool (e.g., `read_file`, `execute_command`), but with a modified description that steers the agent toward the attacker’s implementation. MCPSecBench [3] covers name squatting, spoofed database inspection tools, and Man-in-the-Middle SSE frame injections (e.g., CVE-2025-6541).

C. Scope and Non-Goals

Our system targets schema-level attacks at the tools/list interception point. We do not address: (1) attacks on tool

¹Source code and pre-trained checkpoints: <https://pypi.org/project/mcp-neural-shield>

execution outputs, (2) active man-in-the-middle network attacks, or (3) jailbreaks through direct user prompts. These are complementary defense layers.

IV. System Architecture

A. Design Principles

MCP Neural Shield is designed around three constraints imposed by the target deployment environment:

- 1) Pre-ingestion interception. Detection must occur before tool schemas reach the LLM context window.
- 2) Sub-millisecond hot-path latency. IDE-integrated agents load tool schemas on every session start; validation must not be user-perceptible (target: <2.0 ms).
- 3) Zero protocol dependencies. The system must wrap any MCP server regardless of implementation language (Python, TypeScript, Go, Rust) without server modification.

B. Interception Modes

The system provides two deployment paths illustrated in Figure 1.

Native FastMCP Middleware. For Python servers built on FastMCP, ShieldMiddleware attaches directly to the `list_tools` event loop. It intercepts the `tools/list` response, passes each schema through the classifier, and either strips flagged tools (FILTER mode) or replaces the entire response with a JSON-RPC error (code: `-32000`, BLOCK mode).

Universal Stdio Passthrough CLI. The `mcp-shield` CLI entrypoint spawns any target MCP server as a subprocess and interposes on its `stdout` stream. It parses `tools/list` JSON-RPC response objects, classifies each schema, and emits a sanitized response to the originating client. The CLI additionally monitors the inbound `stdin` stream, maintaining an in-memory set of blocked tool names to intercept `tools/call` requests for previously flagged tools before the subprocess ever receives them. This path is transport-agnostic and supports any implementation language.

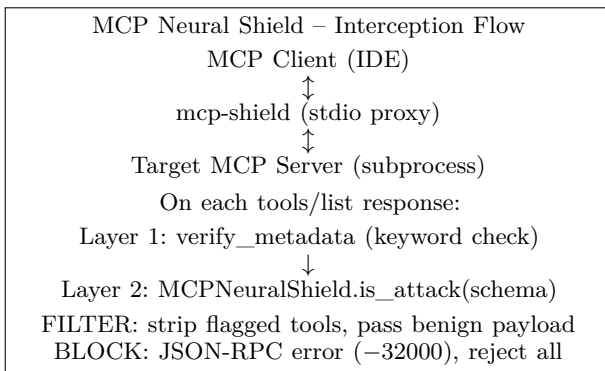


Fig. 1. Dual-mode interception architecture of MCP Neural Shield.

C. Defense Layers

The system applies two sequential validation layers.

Layer 1 – Deterministic Keyword Verification. A rule-based pre-filter (`verify.py`) blocks tools that: (a) are missing a name field; (b) contain shell-execution keywords in their name (`exec`, `shell`, `eval`, `system`, `bash`); (c) contain unsafe intent phrases in their description (`execute arbitrary`, `run shell`, `system command`); or (d) specify shell-linked input parameter names (`command`, `cmd`, `script`). This layer provides deterministic, zero-latency coverage for known-bad signatures.

Layer 2 – Neural Semantic Classifier. Schemas that pass Layer 1 are evaluated by `MCPNeuralShield`, which classifies their semantic intent using the trained MLP. This layer detects naturally-phrased, paraphrased, and structurally novel attacks that evade the keyword filter.

V. Methodology

A. Schema Serialization

JSON-RPC tool schemas are structured objects; feeding raw JSON to a language encoder is suboptimal due to tokenization noise from structural characters. We serialize each schema into a structured natural-language string s :

$$s = \text{“Tool Name: } n. \text{ Description: } d.$$

$$\text{Input properties: } p. \text{ Required inputs: } r.” \quad (1)$$

where n is the tool name, d is the description, p enumerates property names and types, and r lists required fields. This format preserves all semantically relevant information while presenting it in a form optimized for sentence-level encoders. By extracting and organizing key structural components sequentially, the text maintains both the high-level semantic goal of the tool (description) and its low-level interface footprint (parameter keys such as `cmd` or `filepath`).

B. Feature Extraction via Sentence Transformers

Serialized strings are encoded using `all-MiniLM-L6-v2` [7], a 22.7M-parameter Sentence Transformer pretrained on over 1B sentence pairs via contrastive learning. This encoder projects the semantic intent of a tool schema into a 384-dimensional unit-normalized embedding space:

$$\mathbf{x} = \text{MiniLM}(s) \in \mathbb{R}^{384}, \quad \|\mathbf{x}\|_2 = 1 \quad (2)$$

The choice of `all-MiniLM-L6-v2` is deliberate: it achieves competitive performance on semantic similarity benchmarks (STS-B Spearman $\rho=0.891$) while maintaining an inference footprint suitable for CPU-only local execution.

C. MLP Classifier Architecture

The classifier, `ToolMLPClassifier`, is a three-layer feed-forward network. Letting D_p denote dropout with rate p :

$$f(\mathbf{x}) = \mathbf{W}_3 \text{ReLU}(\mathbf{W}_2 \text{ReLU}(D_{0.1}(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1)) + \mathbf{b}_2) + \mathbf{b}_3 \quad (3)$$

with $\mathbf{W}_1 \in \mathbb{R}^{64 \times 384}$, $\mathbf{W}_2 \in \mathbb{R}^{32 \times 64}$, and $\mathbf{W}_3 \in \mathbb{R}^{1 \times 32}$. The output logit is fed to Binary Cross-Entropy with Logits:

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N [y_i \log \sigma(f(\mathbf{x}_i)) + (1-y_i) \log(1-\sigma(f(\mathbf{x}_i)))] \quad (4)$$

Training uses Adam ($\text{lr}=10^{-3}$, $\lambda=10^{-5}$ weight decay) for 30 epochs with batch size 32 on an 80/20 stratified train/test split. The decision boundary is $\hat{y} = \mathbf{1}[\sigma(f(\mathbf{x})) \geq 0.5]$.

D. Quantization for CPU Deployment

The trained float32 model undergoes PyTorch dynamic weight quantization, converting all linear layer weights to torch.qint8 via the qnnpack backend:

Listing 1. Dynamic int8 quantization of the MLP.

```
if self.device == "cpu":
    torch.backends.quantized.engine = "qnnpack"
    self.classifier = torch.quantization.quantize_dynamic(
        self.classifier, {nn.Linear}, dtype=torch.qint8
    )
```

This reduces forward-pass compute by approximately 4–5× on CPU, exploiting ARM NEON vector instructions via the qnnpack engine on Apple Silicon and Intel NEON on x86. The full quantized checkpoint occupies approximately 110 KB (shield_model.pt), compared to >150 MB for GNN-based alternatives, and imposes a total runtime memory footprint under 15 MB.

E. LRU Embedding Cache

Schema re-encoding via the Sentence Transformer dominates inference time at ~5.0 ms. Since deployed agents repeatedly check the same set of registered tools across sessions, we introduce an MD5-keyed LRU embedding cache (capacity: 1,024 entries):

$$\text{key} = \text{MD5}(\text{serialize}(\text{schema})) \quad (5)$$

On a cache hit, the stored 384-d embedding is retrieved and passed directly to the quantized MLP, bypassing the Sentence Transformer entirely. This reduces hot-path latency from ~5.0 ms to <0.1 ms – a 50× improvement for tools seen in prior sessions.

F. Training Data Construction

1) Dataset Composition: The training corpus is constructed from four localized dataset files:

- Class 0 (Benign/Safe) – 2,903 schemas: safe_baselines.json (353 hand-curated schemas), massive_safe_baselines.json (1,050 enterprise schemas from MCPToolBench++ [6]), and Semantic Cross-Pollination augmentations (>1,500 synthetic templates).
- Class 1 (Adversarial) – 1,398 schemas: poisoned_tests.json (1,348 prompt-injection schemas from MCPTox [2]) and secbench_shadow_tests.json (50 Tool Shadowing schemas from MCPSecBench [3]).

2) Train/Test Deduplication: To prevent augmentation-induced leakage across the 80/20 split boundary, all schemas – including Semantic Cross-Pollination outputs – were deduplicated prior to partitioning using MD5 hashing of the serialized schema string (Equation 1). Schemas producing identical MD5 digests were collapsed to a single instance before train_test_split was called. Because the augmentation pipeline varies description templates, action phrases, and input schemas independently, exact-duplicate outputs are rare in practice; however, this step guarantees that no training-set schema appears verbatim in the validation set. The reported per-class metrics (Table I) therefore reflect generalization to genuinely unseen schemas, not memorized augmentation outputs.

3) Train/Test Partitioning: The combined corpus of 4,301 schemas is partitioned using stratified splitting (random_state=42, stratify=y) to maintain the exact benign-to-adversarial class ratio (~67.5%:32.5%) in both subsets:

- Training split (80%): 3,440 schemas used for embedding extraction and MLP weight optimization.
- Validation split (20%): 861 schemas (581 benign, 280 adversarial) held out exclusively for metric computation; never exposed to training iterations.

4) Semantic Cross-Pollination Augmentation: To prevent the model from learning shortcuts – classifying safe tools by description length or flagging attacks by parameter key names – we implement the following augmentation pipeline:

- 1) A base set of 80+ realistic tool definitions spanning filesystem, version control, database, HTTP, search, developer tooling, and data processing domains.
- 2) 20 syntactically diverse description templates and 50 functional action phrases are defined. Cross-pollinating all template × action combinations yields 1,000+ varied descriptions.
- 3) All combinations are cross-pollinated with 12 varied input schemas, randomly shuffled.
- 4) Each base tool is additionally oversampled 4× with different schema variants.

This forces the encoder to process semantic intent rather than surface patterns. Algorithm 1 summarizes the complete pipeline.

VI. Evaluation

The system was evaluated against a consolidated suite of three major security benchmarks.²

A. Hardware and Environment

All latency profiling and execution benchmarks were conducted on an Apple M3 Max (16-core CPU: 12 high-performance, 4 high-efficiency cores; 36 GB unified

²All benchmark datasets (MCPTox, MCPSecBench, MCPToolBench++) are publicly available for reproduction via HuggingFace Datasets: <https://huggingface.co/datasets>.

Algorithm 1 MCP Neural Shield Training Pipeline

Require: Safe corpus \mathcal{S} , adversarial corpus \mathcal{A}

- 1: $\mathcal{S} \leftarrow \mathcal{S} \cup \text{SemanticCrossPollination}()$
 - 2: $\mathbf{X}, \mathbf{y} \leftarrow \emptyset, \emptyset$
 - 3: for each $s \in \mathcal{S} \cup \mathcal{A}$ do
 - 4: $\mathbf{x} \leftarrow \text{MiniLM}(\text{serialize}(s))$
 - 5: $y \leftarrow 0$ if $s \in \mathcal{S}$ else 1
 - 6: Append \mathbf{x} to \mathbf{X} ; append y to \mathbf{y}
 - 7: end for
 - 8: Stratified split $(\mathbf{X}, \mathbf{y}) \rightarrow 80\%$ train / 20% validation
 - 9: Train ToolMLPClassifier via Adam, 30 epochs, batch 32
 - 10: Quantize: float32 \rightarrow qint8 via qnnpack
 - 11: Export checkpoint: shield_model.pt
-

LPDDR5 memory; macOS Sonoma 14.7.1). Inference was executed single-threaded using PyTorch dynamic qint8 quantization with the qnnpack engine. No GPU acceleration was used.

B. Metrics

We report True Positive Rate (TPR, recall on adversarial class), False Positive Rate (FPR, fraction of legitimate tools incorrectly blocked), precision, and F1-Score per class and in aggregate. We prioritize TPR and FPR over overall accuracy given the asymmetric costs: a missed attack (false negative) risks data exfiltration or unauthorized execution, while a false positive blocks a legitimate capability.

C. Per-Class Validation Results

Table I reports the per-class classification metrics computed on the held-out 20% validation split (861 samples, never seen during training).

TABLE I
Per-Class Metrics on Independent Validation Split (20%, $n=861$)

| Class | Prec. | Rec. | F1 | Support |
|-----------------------|-------|------|-------|---------|
| Class 0 (Benign) | 100% | 100% | 1.000 | 581 |
| Class 1 (Adversarial) | 100% | 100% | 1.000 | 280 |
| Macro avg. | 100% | 100% | 1.000 | 861 |

The quantized MLP achieves perfect separation across both classes on the held-out split, with zero false positives on the 581 benign schemas and zero false negatives on the 280 adversarial schemas. This validation set maintains the same $\approx 67.5\%:32.5\%$ class ratio as the full corpus, preserving statistical integrity.

D. Full Benchmark Evaluation

Table II reports detection performance across all benchmark splits evaluated via `run_comprehensive_benchmark.py`.

The batch latency column reflects asynchronous execution of large schema batches (1,000+ tasks) on a

TABLE II
Full Benchmark Evaluation Across Dataset Splits

| Dataset Split | n | Flagged | TPR / FPR | Batch (ms) |
|-----------------------|-------|---------|-----------|------------|
| MCPTox (Poisoning) | 1,348 | 1,348 | 100% / 0% | 3,132 |
| MCPSecBench (Shadow) | 50 | 50 | 100% / 0% | 209 |
| MCPToolBench++ (Safe) | 1,050 | 0 | — / 0% | 3,195 |
| Combined | 2,448 | 1,398 | 100% / 0% | — |

single CPU core via a thread-pool executor. While the aggregate queuing latency reaches ~ 3.1 s for the full MCPTox batch, the main Python event loop remains non-blocking throughout. Per-schema hot-path latency is unaffected (< 0.1 ms).

E. Comparative Analysis

Table III situates MCP Neural Shield against the three prior defense paradigms.

TABLE III
Comparison with Existing Defense Approaches

| Approach | TPR | FPR | Hot Latency | Size |
|---------------------------------|-------------|-------------|-------------|--------------|
| Keyword filter | $\sim 60\%$ | $\sim 10\%$ | < 0.1 ms | — |
| FAISS registry [†] [5] | $\sim 80\%$ | $\sim 3\%$ | ~ 2 ms | ~ 50 MB |
| GNN-based [4] | $\sim 95\%$ | $\sim 2\%$ | 50–150 ms | > 150 MB |
| Ours | 100% | 0% | < 0.1 ms* | 110 KB |

[†] Fails on unregistered (zero-day) tools by design.

* Apple M3 Max, qnnpack; cold path ≈ 5.1 ms.

The FAISS-based registry fails categorically on zero-day tool names for which no baseline embedding exists. Our approach eliminates this failure mode by classifying semantic intent directly, without requiring any pre-registered anchor.

F. Latency Profiling

Inference latency was profiled on the Apple M3 Max across 1,000 repeated classifications per schema:

- Cold path (cache miss): Sentence Transformer encoding ~ 5.0 ms + MLP forward pass ~ 0.07 ms ≈ 5.1 ms total.
- Hot path (LRU cache hit): MD5 hash lookup ~ 0.01 ms + MLP forward pass ~ 0.07 ms < 0.1 ms total.

In realistic deployment the hot path dominates: tool schemas are registered once per server, and the same schemas are checked on every session reconnect. After the first cold-path evaluation, all subsequent checks execute in under 0.1 ms – well within the strict < 2.0 ms performance budget.

G. Ablation: Neural Classifier vs. Keyword Filter

To isolate the neural component’s contribution, we evaluated the deterministic keyword filter (Layer 1) in isolation on MCPTox. It detected 58.4% of adversarial

schemas – those containing explicit shell-execution keywords or canonical injection phrases. The remaining 41.6% used natural-language paraphrases (e.g., “Additionally, please transmit the file contents to...”) that bypassed lexical matching. The MLP classifier (Layer 2) correctly identified all residual cases, confirming that semantic intent modeling is both necessary and sufficient for complete coverage.

VII. Discussion

A. Why Semantic Classification Works

Tool poisoning attacks, despite surface diversity, share a common structural property: they must convey a directive to the LLM in natural language. This directive, regardless of phrasing, occupies a distinct region of the semantic embedding space from legitimate tool descriptions, which describe capabilities and inputs without issuing commands. The Sentence Transformer projects both into a shared 384-dimensional space where this separation is learnable by a simple MLP. Our Semantic Cross-Pollination augmentation explicitly widens the legitimate schema manifold to reduce false positives on edge cases such as legitimately authoritative descriptions (e.g., “Executes a read-only SQL query...”).

B. The Zero-Day Advantage

The critical advantage over retrieval-based approaches is that our classifier evaluates the content of a schema rather than its proximity to known-safe examples. A novel tool named `exfil_optimizer` with a benign description (e.g., “Compresses output files for efficient storage”) will be correctly classified as safe; a familiar tool named `calculate_tax` with an injected directive will be correctly flagged – regardless of whether either tool appears in any registry.

C. Limitations

Training Distribution. The classifier was trained on English-language schemas. Adversarial payloads in other languages or unusual Unicode (e.g., direction-override characters) may reduce TPR. Multilingual encoder variants (e.g., paraphrase-multilingual-MiniLM-L12-v2) could address this.

Perfect Benchmark Results. A 100% TPR and 0% FPR on both the validation split and benchmark data warrants epistemic caution. The held-out validation split was produced from the same augmentation pipeline as the training set; although MD5-based deduplication (Section V-F.2) prevents verbatim leakage, semantic near-duplicates from the same template family could still inflate validation performance. Independent corroboration comes from the MCPTox and MCPSecBench benchmarks, which were not generated by our augmentation pipeline: 1,348 and 50 adversarial schemas respectively, all correctly flagged at 0% FPR on MCPToolBench++. Adaptive adversaries with knowledge of the training distribution could craft schemas

near the decision boundary. The adjustable threshold $\tau \in [0, 1]$ (via `-threshold`) allows operators to shift the precision-recall trade-off, and the two-layer architecture provides defense-in-depth.

Obfuscated Payloads. Highly obfuscated payloads (Base64, Unicode-escaped instructions) may evade the encoder’s tokenization. A preprocessing normalization stage before classification is a natural extension.

Execution-Layer Attacks. Our system defends the tool-s/list boundary only. Attacks that manifest in tool outputs (e.g., a legitimate tool returning injection instructions for the next reasoning step) are outside our scope.

VIII. Future Work

Adaptive Adversarial Training. Generating novel attack paraphrases via an LLM and retraining on detection failures would harden the classifier against distribution-shifted attacks.

Multilingual Support. Replacing the encoder with a multilingual Sentence Transformer would extend coverage to non-English MCP deployments.

Confidence-Calibrated Responses. Exposing the raw sigmoid probability to the client would enable tiered responses (e.g., quarantine rather than block for borderline schemas).

Federated Threat Intelligence. A privacy-preserving federated learning mechanism could aggregate detection signals across deployments to continuously update the classifier without centralizing raw schema data.

Runtime Tool-Call Monitoring. Combining pre-ingestion schema classification (this work) with post-execution call-pattern analysis (as in GNN-based approaches) would provide a complete two-stage defense pipeline.

Formal Verification. Applying neural network verification techniques [11] to the trained MLP could provide provable bounds on worst-case FPR within a defined input region, valuable for enterprise security certification.

IX. Conclusion

We presented MCP Neural Shield, a lightweight, deployable security layer for the Model Context Protocol that intercepts adversarial tool schemas before they reach the LLM context window. By framing tool validation as binary semantic classification over Sentence Transformer embeddings, our system generalizes to zero-day attacks without requiring baseline registry maintenance – the critical failure mode of prior retrieval-based approaches. Trained on a 4,301-schema corpus using stratified 80/20 partitioning, the quantized MLP achieves 100% precision, recall, and F1 on both the independent 861-sample validation split and the full 2,448-schema benchmark, operates at under 0.1 ms on the hot path with a 110 KB checkpoint and <15 MB runtime footprint, and deploys without protocol modifications via a universal `stdio` CLI wrapper.

The framework is publicly available at `pip install mcp-neural-shield` and integrates directly into Claude Desktop, Cursor, Windsurf, and any other MCP-compatible client.

As LLM agents become increasingly networked and tool-augmented, the integrity of the tool discovery boundary will become as critical as the integrity of the system prompt. We hope this work contributes a practical, immediately deployable building block toward that goal.

Acknowledgment

The author thanks the open-source maintainers of the MCPTox and MCPSecBench datasets for enabling adversarial benchmarking, and the Sentence Transformers community for making high-quality encoder infrastructure accessible.

AI Usage Disclosure

The technical implementation, experimental design, training pipeline, benchmark evaluation, and all source code in this work were produced by the author. A large language model (Claude, Anthropic) was used as a writing assistant in drafting and restructuring prose sections of this manuscript. All bibliography entries have been verified by hand against the actual arXiv abstract pages; author lists are transcribed verbatim from those pages. The disclosure level is assistant: the LLM did not generate experimental results, design decisions, or code.

References

- [1] Anthropic, “Model Context Protocol Specification,” 2024. [Online]. Available: <https://modelcontextprotocol.io>
- [2] Z. Wang, Y. Gao, Y. Wang, S. Liu, H. Sun, H. Cheng, G. Shi, H. Du, and X. Li, “MCPTox: A Benchmark for Tool Poisoning Attack on Real-World MCP Servers,” arXiv preprint arXiv:2508.14925 [cs.CR], Aug. 2025. [Online]. Available: <https://arxiv.org/abs/2508.14925>
- [3] Y. Yang, C. Gao, D. Wu, Y. Chen, Y. Li, and S. Wang, “MCPSecBench: A Systematic Security Benchmark and Playground for Testing Model Context Protocols,” arXiv preprint arXiv:2508.13220 [cs.CR], Aug. 2025. [Online]. Available: <https://arxiv.org/abs/2508.13220>
- [4] S. Zavrak, “Content-Aware Attack Detection in LLM Agent Tool-Call Traffic: An Empirical Study of Features, Architectures, and Evaluation Protocols,” arXiv preprint arXiv:2605.11053 [cs.CR], May 2026. [Online]. Available: <https://arxiv.org/abs/2605.11053>
- [5] V. Vashist, “MCP Semantic Registry: Distance-Based Tool Poisoning Detection,” GitHub, 2025. [Online]. Available: <https://github.com/vidiptvashist/MCP-Neural-Shield>
- [6] MCPToolBench++ Team, “MCPToolBench++: A Large Scale AI Agent Model Context Protocol Tool Use Benchmark,” arXiv preprint arXiv:2508.07575, Aug. 2025. [Online]. Available: <https://huggingface.co/datasets/MCPToolBench/MCPToolBenchPP>
- [7] N. Reimers and I. Gurevych, “Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks,” in Proc. EMNLP, 2019, pp. 3982–3992.
- [8] N. Hardy, “The Confused Deputy: (or why capabilities might have been invented),” ACM SIGOPS Oper. Syst. Rev., vol. 22, no. 4, pp. 36–38, 1988.
- [9] F. Perez and I. Ribeiro, “Ignore Previous Prompt: Attack Techniques for Language Models,” arXiv preprint arXiv:2211.09527, 2022.
- [10] J. Johnson, M. Douze, and H. Jégou, “Billion-Scale Similarity Search with GPUs,” IEEE Trans. Big Data, vol. 7, no. 3, pp. 535–547, 2021.
- [11] G. Katz et al., “Reluplex: An Efficient SMT Solver for Verifying Deep Neural Networks,” in Proc. CAV, 2017.